

1 We thank the reviewers for their constructive feedback. R1 and R2 appreciated the novelty of our unified framework that  
 2 solves both feature imputation and label prediction aspects for the missing data problem, and all reviewers agreed that  
 3 our solution is well motivated. Reviewers pointed out several main concerns, which we summarize and answer below:

4 **1 Novelty of the GNN method (R3 R4).** We thank R3 and R4 for noticing that the core GNN components have been  
 5 separately used in other applications, and the formulation as a bipartite graph has been used for matrix completion task.  
 6 However, we emphasize that our main contribution is *not the particular GNN model but the graph-based framework*.  
 7 We show that a seemingly unrelated missing data problem (imputation and learning subsequent tasks) can naturally  
 8 be solved with graphs and we propose the first graph-based solution to it, which is acknowledged by R1 and R2.

9 Nevertheless, we agree with R3 and R4 that the justification  
 10 of the adopted method can be further improved and we will  
 11 do that in the final version of the paper. **New results to**  
 12 **justify our model.** In addition to ablation study in Section  
 13 4.6, we further justify our model by showing how different types of aggregation, *i.e.*, MEAN(GraphSAGE-mean),  
 14 SUM(GIN), MAX(GraphSAGE-pool) affect the performance. These new results justify our selection of GraphSAGE-  
 15 style pooling in our architecture. **Justifying the new results.** While SUM is theoretically most expressive, in our setting  
 16 the degree of a specific node is determined by the number of missing values which is random and unrelated to the  
 17 missing data task; in contrast, the MEAN and MAX aggregators are not affected by this inherent randomness of node  
 18 degree. We will add these studies to justify our framework.

19 **2 Comparing with more baselines (R1 R2 R3).** We thank reviewers for pointing out additional baselines to compare,  
 20 in addition to the 7 baselines that we have already examined (lines 169–181). **New baseline results.** Summarizing  
 21 reviewer’s suggestions, we additionally compared our methods with 3 representative baselines on the feature imputation  
 22 task: (1) low-rank matrix completion method for mixed-type data (missMDA), (2) deep latent variable models and  
 23 autoencoders (MIWAE), (3) GNN-based baselines (GC-MC (Berg et al 2017) [1], IGMC (Zhang et al 2019) [2]).  
 24 **Discussion.** Note that *all the new baselines are only applicable to the feature imputation task* not the label prediction  
 25 task, and *the GNN-based baselines are only applicable to discrete-value data*, as they explicitly use different weights  
 26 for each distinct value, and thus cannot apply to mix-typed data; GRAPE *has none of these limitations*. We find that our  
 27 GRAPE framework can consistently outperform these new baselines in all the datasets. In 3 discrete-value datasets,  
 28 IGMC [2] has advantages as it is specifically designed for discrete-value data, while our model can still outperform  
 GC-MC [1]. We will include these new results and make sure to cite all the papers suggested in the reviews.

	concrete	energy	housing	kin8nm	naval	power	protein	wine	yacht
Sum	0.094	0.143	0.078	0.277	0.024	0.134	0.040	0.069	0.154
Max	<b>0.088</b>	0.142	<b>0.074</b>	0.252	<b>0.006</b>	<b>0.102</b>	<b>0.024</b>	<b>0.063</b>	0.153
Mean	0.090	<b>0.136</b>	0.075	<b>0.249</b>	0.008	<b>0.102</b>	0.027	<b>0.063</b>	<b>0.151</b>

	concrete	energy	housing	kin8nm	naval	power	protein	wine	yacht	Flixster	Douban	Yahoo	
missMDA	0.190	0.225	0.142	0.285	0.038	0.215	0.068	0.090	0.226	GC-MC [1]	0.917	0.734	20.5
MIWAE	0.156	0.153	0.098	0.262	0.020	0.117	0.042	0.087	0.224	IGMC [2]	<b>0.872</b>	<b>0.721</b>	<b>19.1</b>
Ours	<b>0.090</b>	<b>0.136</b>	<b>0.075</b>	<b>0.249</b>	<b>0.008</b>	<b>0.102</b>	<b>0.027</b>	<b>0.063</b>	<b>0.151</b>	Ours	<b>0.899</b>	<b>0.733</b>	<b>19.4</b>

29 **3 Scalability (R2 R3 R4).** We thank reviewers for asking to discuss the scalability of our method. We pick UCI  
 30 as they are widely-used datasets for benchmarking imputation methods, with *both discrete and continuous features*.  
 31 Our GRAPE *framework does not suffer from scalability issues*: the core of our method is a GNN, which has been  
 32 successfully applied to graphs with billions of edge; while R3 pointed out that a naive one-hot scheme will not work  
 33 when the number of features is prohibitively large, we can use an embedding matrix to encode each feature, which has  
 34 successfully handled millions of tokens in NLP. We agree with reviewers that experimenting on datasets with more  
 35 features can justify the scalability. We provide new results on large-scale benchmarks (Flixster: 2956 features, Douban:  
 36 3000 features, Yahoo: 1363 features) shown in the table above. Note that these datasets *only have discrete features*, and  
 37 therefore they were not used in our manuscript. We will include these new discussions and results in the revised paper.  
 38

39 **4 Discussions of related work (R1).** We thank R1 for pointing out that our discussions for the related work can be  
 40 improved. We intended to give an overview for common limitations in lines 23–31 by starting with “often exhibit  
 41 notable shortcomings”, but we do agree that this summary should be more rigorous. **Clarifications.** (1) We will cite  
 42 these matrix completion models designed for both discrete and continuous variables, and for online learning. (2) We  
 43 will change to “imputation approaches based on deep generative models do not *explicitly* use feature values from other  
 44 observations”. (3) After formulating the missing data problem as a bipartite graph over feature and observation nodes,  
 45 the proposed augmented node feature initialization is natural: feature nodes are not permutation invariant so one-hot  
 46 tokens are used; observation nodes are permutation invariant so constant node features are used – we have provided a  
 47 thorough discussion in lines 130–149. Overall, we will include these discussions in the revised paper.

48 **5 Clarifications.** **Q: (R2)** “Applicable to vision data” **A:** Yes, *e.g.*, our model can be jointly used with CNNs and  
 49 trained end-to-end. **Q: (R2 R4)** “Theoretical guarantees” **A:** The success of our framework is supported by the universal  
 50 approximation capabilities of GNNs on graph structured data. **Q: (R4)** “Random missing data assumption” **A:** This is  
 51 the most common evaluation regime used in missing data papers. Our model can apply to other missing data scenarios.  
 52 **Q: (R4)** “Using labels in feature imputation” **A:** We agree that label information could be helpful in feature imputation  
 53 but we do not use it in our experiments, since the common evaluation for feature imputation tasks does not use labels.